



TITLE:

処理効率向上のためのネットワーク構造の変換(知識ベースとデータベースの統合化に関する研究)

AUTHOR(S):

古川, 哲也; 斎藤, 邦子; 上林, 彌彦

CITATION:

古川, 哲也 ...[et al]. 処理効率向上のためのネットワーク構造の変換(知識ベースとデータベースの統合化に関する研究). 数理解析研究所講究録 1986, 593: 18-28

ISSUE DATE:

1986-06

URL:

<http://hdl.handle.net/2433/99513>

RIGHT:

処理効率向上のためのネットワーク構造の変換

九州大学工学部 古川哲也 (Tetsuya Furukawa)

九州大学工学部 斎藤邦子 (Kuniko Saito)

九州大学工学部 上林彌彦 (Yahiko Kambayashi)

1. まえがき

データベースにおけるデータモデルとしては、関係モデル、ネットワークモデルが代表的であるが、処理効率は一般にネットワークモデルの方がよく、ネットワークデータベースシステムは実用的システムとして広く普及している。しかしネットワークモデルは構造によっては処理効率が大きく異なり、本稿ではネットワークモデルの構造と質問処理の適合性を明確にし、与えられた質問を効率よく処理できるようにネットワーク構造を変換する方法を示す。

データのモデル化での重要な問題として、意味制約の保持と質問処理の効率化があげられる。関係モデルでは2つの問題を分離しており、意味制約はデータベース設計時に構造に反映され、質問処理の効率化は質問最適化によっている。一方ネットワークモデルでは、構造に制約があり、処理効率も巡航操作による処理のため構造によって異なる。従って、ネットワーク構造の決定（スキーマ設計）ではこれら両方の問題を考慮する必要がある。スキーマ設計では、関数従属性集合や多値従属性集合による設計^{[1][2][3]}が研究されているが、処理効率まで考慮したものは知られていない。

本稿では、与えられた質問を効率よく処理できるようにするためのネットワーク構造の変換について、意味制約（従属性制約）の保持に関するデータベースの管理の問題を含めて議論する。処理効率の向上は主に冗長性の付加によるが、冗長性に対しては管理を必要とする。本稿での変換法は、効率よい処理を行なうためにネットワーク構造が満たすべき条件を与え、質問処理の効率化と冗長性に対する管理のトレードオフの問題から変換後のネットワーク構造がどの条件を満たすようにするかを選択できるようにしている。

2節で例題を用いて構造による処理効率の違いを示して問題を明確にし、3節でネットワークモデルに関する基本的事項を示す。4節ではネットワーク構造と質問処理の適合性について議論し、効率よい質問処理のためのネットワーク構造が満たすべき条件を与える。5節でネットワーク構造の変換法を示し、6節では変換後のネットワーク構造が持つ性質について検討する。最後に7節でまとめを行なう。

2. 質問処理に適したネットワーク構造

属性集合ABCD に対しAB の値を指定して対応するCDの値を求める質問 ($Q(AB, CD)$ で表わす) を考える。この質問の処理で最も効率がよいと思われるのは、図1(a)のネットワーク構造である。ここで、節点 $R(X)$ は属性集合Xからなるレコードの集合であり、有向枝はその方向に

1対多のレコードの対応があることを表わす。また下線はキー属性である(詳細な定義は3節で行なう)。この構造では、 R_1 でA の値によりレコードが唯一決定まり、 R_2 でそのレコードに対応しB の値の指定を満たすレコードが定まる。それに対応する R_3 のレコード集合のCDの値がこの質問の解となる。

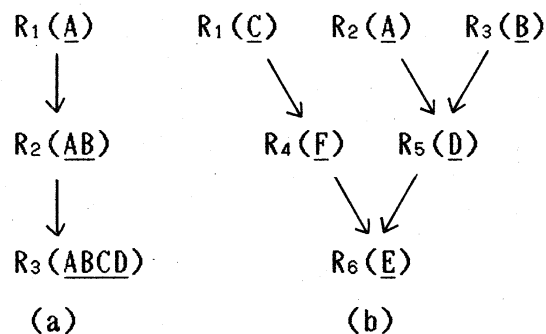


図1 基本的なネットワーク構造

属性の全体集合を $U=ABCDEF$, U で満足される関数従属性の集合を $\{E \rightarrow DF, F \rightarrow C, D \rightarrow AB\}$ とすると、この集合を構造に反映させたものは図1(b)となる。この構造では、質問 $Q(AB, CD)$ の処理効率はあまりよくない。その理由として次のものがある。

- (1) AB の値の指定を満たすものは、 R_2, R_5, R_3 の経路で求めることができるが、図1(a)では R_1, R_2 で得られたのに対し検索するレコード型が多い。
- (2) $R_4(F)$ や $R_6(E)$ は解には関係ないが、レコードの対応を求めるために検索する必要がある。
- (3) $R_1(C)$ と $R_5(D)$ のレコードの対応は R_6 により多対多となり、解が重複する。

図1(a),(b)を併合して図2(a)とし、質問処理は R_2, R_7, R_8 で、意味制約の保持

は $R_1 \sim R_6$ で行なうようにするとすることもできる。ここで、枝 $(R_1, R_8), (R_5, R_8)$ は一貫性を保つために必要なものである。このときは次の問題がある。

- ・ R_8 のCDの値の対応は、 R_1, R_4, R_6, R_5 の経路で得られるCDの対応と一致しなければならない。
- ・ R_7, R_8 でのABCDの値の対応は、 R_2, R_3, R_5, R_8 の経路でのABCDの対応と一致しなければならない。

即ち、閉路が

できると管理の問題が生じる^[4]。

本稿の方法では、ネットワーク構造は

図2(b)に変

換される。こ

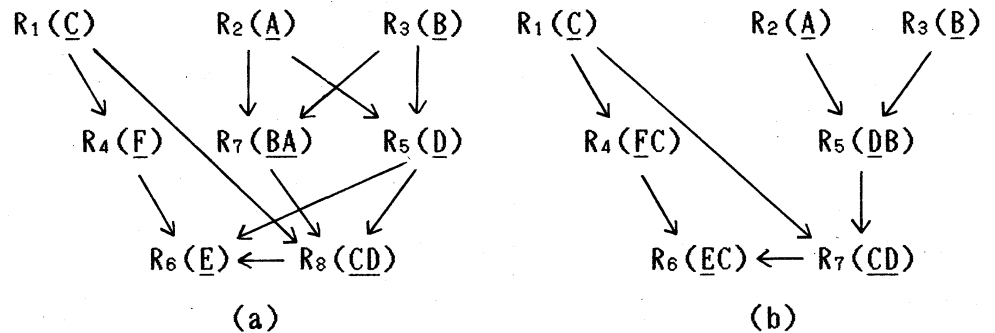


図2 質問処理と意味制約保持のための構造

の構造では、表現する従属性は図1(b)に等しい。また、質問処理は R_2, R_5, R_7 の経路で行なうことができ、図1(b)のような問題は生じない。このときも閉路ができるが、各枝のレコードの対応(親子集合)で同じ属性の値(ここではCの値)が同じとなるようにすることによって管理できる^[4]。即ち、意味制約を表わす構造をなるべく残して質問処理の効率化をはかるものである。

3. 基本的事項

ネットワークモデルは、同じデータ項目(本稿では関係モデルとの対応のため属性と呼ぶ)からなるレコードの集合であるレコード型と、2つのレコード型間のレコードの1対多の対応を表す親子集合型の集合によって表現される。属性集合 X で構成されるレコード型 R を $R(X)$ で、親レコード型が R_0 、子レコード型が R_m である親子集合型を $\langle R_0, R_m \rangle$ で表わす。 $\langle R_0, R_m \rangle$ では R_0 の各レコードに対し R_m の任意個のレコードが対応し、それぞれを親子集合という。

ネットワークモデルの構造はバックマン線図と呼ばれる有向グラフ $B(V, A)$ で表わされる。 V は各レコード型に対応する節点集合、 A は各親子集合型に対応し、

親レコード型に対応する節点から子レコード型に対応する節点に向かう有向枝の集合である。値が定まればレコード型 $R(X)$ の対応するレコードがただ1つ定まるような最小の属性集合を、 $R(X)$ のキーと呼び $K(R(X))$ で表わす。必ずしも $K(R(X)) \subseteq X$ である必要はなく、 R が子レコード型となる親子集合が定まれば（親レコード集合が定まれば）その子レコード間で $K(R) \cap X$ の値によってレコードが一意に定まればよい。バックマン線図中ではキーを下線をつけた属性集合（ X に含まれない属性は括弧内）で表わす。単一属性を A, B, \dots で、属性集合を \dots, X, Y, Z で表わし、属性集合の和集合を連接で表わす。

本稿で対象とする質問は、属性集合 X_s の値を指定し対応する属性集合 X_p の値を求めるもので、 $Q(X_s, X_p)$ で表わす。 X_s, X_p はそれぞれ関係代数での選択、射影に対応する。質問 Q は次の定義により明確にされる。

〔定義〕 部分ネットワーク構造

バックマン線図 B の連結な部分グラフで表わされるネットワーク構造を B の部分ネットワーク構造という。特に、属性集合 X に対して X と共通属性を持つレコード型とその間の経路のレコード型、親子集合型からなる部分ネットワーク構造を $B(X)$ で表わす。

〔定義〕 レコードの組、ネットワーク構造における関係

バックマン線図 B で表わされるネットワーク構造に含まれるレコード型を R_1, R_2, \dots, R_n とする。レコード集合 $t = \{r_1, r_2, \dots, r_n\}$ （ r_i は R_i のレコード（ $1 \leq i \leq n$ ））で B に親子集合型 $\langle R_i, R_j \rangle$ があれば、 r_i は r_j の親レコードであるとき、 t を B におけるレコードの組と呼ぶ。また、 B に含まれるレコード型を構成する属性集合の和集合を $U = \{R.A \mid \text{レコード型 } R \text{ の属性 } A\}$ とする。 B の関係は U 上の関係 u であり、 $u(B)$ で表わす。 $u(B)$ の各組は B におけるレコードのすべての組と1対1に対応し、対応する組とレコードの属性の値は等しい。

$Q(X_s, X_p)$ は、解がある部分ネットワーク構造 B の関係 $u(B)$ に対する X_s による選択及び X_p への射影演算で求められるものである。このような質問は、関係モデルにおける連結な自然結合質問のクラスとほぼ対応しており^[6]、一般的なクラ

スである。

〔定義〕 質問スキーマ

質問 Q の解が、部分ネットワーク構造 B の関係 $u(B)$ に対する選択及び射影演算で求められるとき、 B を Q の質問スキーマ B_Q という。

バックマン線図 B ，質問 Q に対し、 B の冗長性により B における Q の質問スキーマは一意ではない。非巡回質問スキーマが存在する質問を木質問，それ以外を巡回質問と呼ぶ^[5]。巡回性はバックマン線図を無向としたときのものであり、本稿では木質問を対象とする。

ネットワーク構造はレコード型と親子集合型の集合であるため、関数従属性集合と多値従属性集合（結合従属性）を表現しており^[4]、本稿では特に関数従属性について議論する。関数従属性は属性集合 X, Y に対し $X \rightarrow Y$ で表わされ、 X の値を決めると対応する Y の値がただ 1 つ定まるというものである。1 つのレコード型 $R(X)$ では関数従属性 $K(R) \rightarrow X$ を表わしており、1 つの親子集合型 $\langle R_0(X_0), R_m(X_m) \rangle$ では $K(R_m) \rightarrow X_0$ を表わしている。従って、レコード型 R に対しその祖先のレコード型に含まれる属性集合を $\text{anc}(R)$ とすると、 $K(R) \rightarrow \text{anc}(R)$ となる。

4. 質問処理とネットワーク構造の適合性

質問 $Q(X_s, X_p)$ の処理効率が悪くなるのは、次の場合がある。これらは 2 節の図 2 (b) での問題点に対応している。

- (1) X_s の条件を満たすかどうかの検査が、その属性を含むレコード型が分散しているために遅くなり、解とはならないレコードの対応を処理中に求めなければならない。
- (2) 巡航中に $X_s X_p$ 以外の属性集合からなるレコード型も巡航する。
- (3) 解が重複する。

(1) は、 B_Q 中で X_s の属性を含むレコード型がそれのみで連結であればよく、(2) は $X_s X_p$ の属性を含むレコード型がそれのみで連結であればよい。従って (1), (2) に対するネットワーク構造が満たすべき条件は次のようになる。

[条件1] バックマン線図Bと質問 $Q(X_S, X_P)$ に対し、 $B_Q(X_S)$ 中の各レコード型 $R(X)$ について、 $X_S \cap X \neq \emptyset$.

[条件2] バックマン線図Bと質問 $Q(X_S, X_P)$ に対し、 B_Q 中のレコード型 $R(X)$ について、 $X_S X_P \cap X \neq \emptyset$.

(3)の解の重複については次の補題と定理による。

[補題1] ネットワーク構造Bにおいて、子レコード型を持たないレコード型を R_1, R_2, \dots, R_n とする。B中の全属性集合を U , $K = \bigcup K(R_i) - \bigcup (\text{anc}(R_i) - K(R_i))$ としたとき、 K は $K \rightarrow U$ となる最小の属性集合である。

(証明) レコード型 R_i の属性集合を X_i とする。

$K(R_i) \rightarrow \text{anc}(R_i)$, $K(R_i) \rightarrow X_i$ ($1 \leq i \leq n$) 及び、 $U = \bigcup (X_i \cup \text{anc}(R_i))$ より、

$\bigcup K(R_i) \rightarrow U$ となる。また、 $K(R_i) \rightarrow K(R_j) \cap \text{anc}(R_i)$ なので、

$\bigcup K(R_i) - \bigcup (\text{anc}(R_i) - K(R_i)) \rightarrow U$ である。

[定義] 補題1中の K をネットワーク構造Bのキーと呼び、 $K(B)$ で表わす。

[定理1] 質問 $Q(X_S, X_P)$ の解がどのようなデータベースの状態と X_S に対しても重複がないための必要十分条件は、 $K(B_Q) \subseteq X_S X_P$ である。

(証明) $Q(X_S, X_P)$ の解に常に重複がなければ、 $X_S X_P$ の値の組は B_Q のレコードの組を一意に定める。即ち、 $X_S X_P$ は B_Q の関係 $u(B_Q)$ の超キーとなるので $K(B_Q) \subseteq X_S X_P$ である。 $K(B_Q) \subseteq X_S X_P$ であれば、 $K(B_Q)$ に重複がないので $X_S X_P$ の値の組にも重複がない。 X_S の値を1つ定めたとき、対応する X_P の値に重複があれば $X_S X_P$ に重複があることになるので、 X_P の値、即ち解に重複はない。

[条件3] バックマン線図Bと質問 $Q(X_S, X_P)$ に対し、 $K(B_Q) \subseteq X_S X_P$.

これらすべての条件を満たすようにすると、冗長性が増し、複雑なネットワーク構造となることがある。質問処理にそれほどの効率を要求しないならば、管理を簡単にするために冗長性は少ない方がよい。このような管理とのトレードオフから、すべての条件を満たすようにネットワーク構造を変換するのではなく、指定された条件のみを満たすように変換する。

5. ネットワーク構造の変換方法

条件 1, 2 に対するネットワーク構造の変換は、共に次の変換 1 を用いて行なうことができる。

[変換 1] バックマン線図 B , 属性集合 Z に対し、次の変換を行なう。仮想質問 $Q(\phi, Z)$ を考える。

0. B_0 中で Z と共通属性を持たないレコード型を R とする。
 1. R の親レコード型, 子レコード型が B_0 中にそれぞれ 1 つずつ (R_0, R_m とする) のとき、 B, B_0 に親子集合型 $\langle R_0, R_m \rangle$ を加え、 B_0 から R を除く。 $\langle R_0, R_m \rangle$ を marked とする。 $\langle R_0, R \rangle, \langle R, R_m \rangle, R$ が marked であれば、それを B から除く。
 2. R の親レコード型が B_0 中になく、子レコード型を R_1, R_2, \dots, R_n (必ず複数ある) とする。 B, B_0 にレコード型 $R'(X')$ ($X' = (\bigcup_i X_i \cap Z) \cup \bigcup_i K(R_i)$) , 親子集合型 $\langle R_1, R' \rangle, \langle R_2, R' \rangle, \dots, \langle R_n, R' \rangle$ を加える。 B_0 から R を除く。 R が marked であれば、それを B から除く。
 3. その他のとき、 R の適当な親レコード型 $R_0(X_0)$ とする。 R に $Z \cap X_0$ を加える。 R_0 が B_0 中で他に連結なレコード型を持たなければ R_0 を B_0 から除く。
- 0 ~ 3 を変換できなくなるまで繰り返す。

条件 1, 2 を満たすようにするには、 B を B_0 , X をそれぞれ $X_S, X_S X_P$ として変換 1 を行なえばよい。条件 2 に対しては、変換 1 中の仮想質問の質問スキーマは B_0 に一致する。

[補題 2] 変換 1 で条件 1, 2 を満たすように変換できる。

(例 1) 図 1(b) のネットワーク構造と質問 $Q(AB, CD)$ に対し、条件 1 を満たすように変換する。仮想質問に対する質問スキーマは R_2, R_3, R_5 のレコード型からなり、 R_5 に R_3 の属性 B を加えることにより条件 1 が満たされる。さらに条件 2 を満たすように変換すると図 3(a) となり、このときの質問スキーマは図 3(b) である。

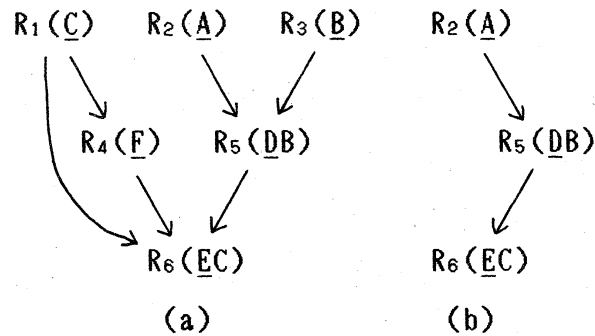


図 3 条件 1, 2 を満たす構造

条件 3 の解の一意性に対しては、次の変換 2 を用いることによって条件を満たすようにすることができる。

〔変換 2〕 バックマン線図 B , 属性集合 Z , 仮想質問 $Q(\phi, Z)$ に対し、 $K(B_0)$ が Z の部分集合となるように変換する。

0. B_0 中で、 $(K(B_0) - Z) \cap K(R) \neq \phi$ となり、子レコード型を持たないレコード型を $R(X)$, R の親レコード型を $R_1(X_1), R_2(X_2), \dots, R_n(X_n)$ とする。 $X_i \cap Z = \phi$ のとき、変換 1 を用いて R_i を除くか Z と共通集合を持つように変換する。

1. レコード型 $R'(X')$ ($X' = (X \cup \bigcup X_i) \cap Z$), 親子集合型 $\langle R', R \rangle$ を作る。

2. $K(R_i) \subseteq Z$ のとき、親子集合型 $\langle R_i, R' \rangle$ を作る。 $\langle R_i, R' \rangle$ の R_i と R のレコードの対応は $\langle R_i, R' \rangle, \langle R', R \rangle$ による R_i と R の対応に一致するので、 $\langle R_i, R \rangle$ を B から除く。

$K(R_i) \not\subseteq Z$ のとき、レコード型 $R_i'(X_i')$ ($X_i' = K(R_i) \cap Z$) と親子集合型 $\langle R_i', R' \rangle, \langle R_i', R' \rangle$ を作る。

0 ~ 2 を繰り返す。

質問 $Q(X_S, X_P)$ に対しては、変換 2 は $B = B_0$, $Z = X_S X_P$ として行なう。

〔補題3〕 変換2で条件3を満たすように変換できる。

(例2) 図3(a)の構造と質問 $Q(AB, CD)$ に対し、条件3を満たす質問スキーマが存在するように変換する。質問スキーマ B_0 をどのようにとっても $K(B_0)$ は R_6 のキー E となる。変換後の構造は図4(a)であり、このときの質問スキーマは図4(b)となる。

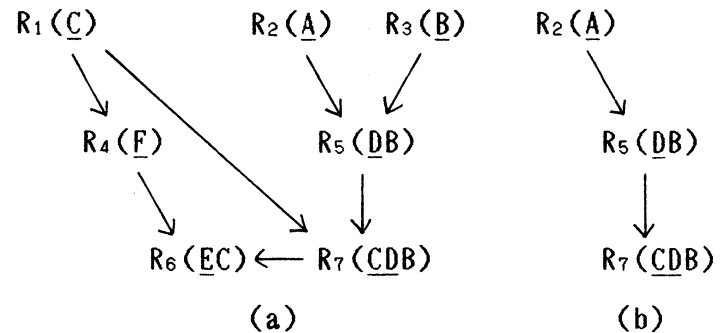


図4 条件3を満たす構造

6. 変換後のネットワーク構造の性質

本節では、5節で示したネットワーク構造の変換に関する性質を、変換の可能性と従属性の保持について議論する。

〔補題4〕 質問 Q に対し、条件1を満足するネットワーク構造を条件2を満たすように変換1を用いて変換した構造は、条件1を満足する。

(証明) 変換1で新たな経路が作られるのは1と2の場合である。両方とも与えられた属性集合を含まないレコード型を除くもので、条件1に対する性質を変えない。

〔補題5〕 質問 Q に対し条件1, 2を満足するネットワーク構造を条件3を満たすように変換2を用いて変換した構造は、条件1, 2を満足する。

(証明) 変換2の1では、条件1又は2を満足する方向の変換である。また、質問スキーマから除かれるレコード型 R に含まれる属性は、その前に加えられたレコード型に含まれ、そのレコード型は R と連結なレコード型と連結になるので、条件1, 2に対する性質を変えない。

〔定理 2〕 任意の質問 Q に対し、バックマン線図 B を条件 1, 2, 3 のうち指定されたものを満たすように変換できる。

(証明) 補題 2, 3 より、各条件に対してはそれを満たすように変換できる。また、補題 4, 5 より、条件に対する変換を 1, 2, 3 の順で行なえば指定された条件を満たすようにすることができる。

以上のように指定された条件を満足するようにネットワーク構造を変換することはできるが、変換後も処理効率があまり変わらない場合がある。変換 1 の 3 で R_0 を B_0 から除けない場合がそれで、このときは条件指定の検査や解を求めるために巡航するレコード型は変わらない。さらに効率よくするために、質問スキーマが(有向の)木構造になるようにすることもできる。これは、変換 1 の 2 と 3 を属性集合 Z と共通属性を持つものに対しても行なうことによる。さらに索引構造の付加による効率化^[4]も行なうことができる。

ネットワーク構造に変換 1, 2 を行なっても、その表現する従属性制約は変わらない。変換 1 ではもとのネットワーク構造から除かれるものはなく、制約の保持をもとの構造で行ない、加えられた構造はもとの構造から得られる対応を表わす。変換 2 では、除かれる親子集合型の表わす関数従属性は加えられた構造で保存され、その構造が表わす制約ももとの構造から得られるものである。

〔定理 3〕 変換 1, 2 はネットワーク構造が表現する従属性制約を変えない。

冗長な部分に対する管理は、データの更新のときに問題になる。特に、バックマン線図に閉路ができたときは複数の経路でレコードの対応が同じになる必要がある。しかし、閉路内の他のレコード型の子レコード型とならないもの 1 つのときは、そのレコード型のキーを閉路のすべてのレコード型に付加することによって管理が簡単になる^[4]。変換 2 でできる閉路は必ずそのようなものとなり、変換 1 でできる閉路も親子集合型の付加によりそのような閉路集合にすることができる。

(例3) 図4(a)では閉路が存在するが、 R_1 のキーCを閉路内のすべてのレコード型に付加することにより管理を簡単にできる。結果は、図2(b)となる。

7. あとがき

質問処理を効率化するために、ネットワーク構造を変換する方法を示した。この変換は、ネットワーク構造が表現する関数従属性を変えないことを示したが、容易に多値従属性、結合従属性に拡張できる。この方法は、ネットワーク構造の設計時だけでなく、データベースの使われ方の変化に対しても適用できる。

変換1, 2は選択するレコード型などが非決定的となっており、これは今後解決すべき問題である。これは質問スキーマ B_0 での巡航順序の決定に左右され、最適な巡航順序との関係を明確にしなければならない。また、本稿で対象とした質問のクラスを拡張し、関係代数で記述される質問のクラスと等価な質問^[6]についても処理効率化のためのネットワーク構造の変換法を明確にしなければならない。

参考文献

- [1] Lien, Y.E., "On the Equivalence of Database Models", J.ACM, vol.29, no.2, pp.333-362, April 1982.
- [2] Kuck, S.M., Sagiv, Y., "Designing Globally Consistent Network Schemes", Proc. ACM SIGMOD Int. Conf. on Management of Data, pp.185-195, May 1983.
- [3] Yannakakis, M., "Algorithms for Acyclic Database Schemes", Proc. 7th Int. Conf. VLDB, pp.82-94, Sept. 1981.
- [4] Kambayashi, Y., Furukawa, T., "Semantic Constraints Expressed by Network Model", Proc. Foundations of Data Organization, pp.201-206, May 1985.
- [5] 古川, 上林, "ネットワークデータベースにおける木質問", 電子通信学会技術研究報告, AL84-62, 1985年1月.
- [6] 古川, 上林, "ネットワークデータベースにおける質問の複雑さについて", 情報処理学会研究報告, DB-52-5, 1986年3月.